

# Object Localization

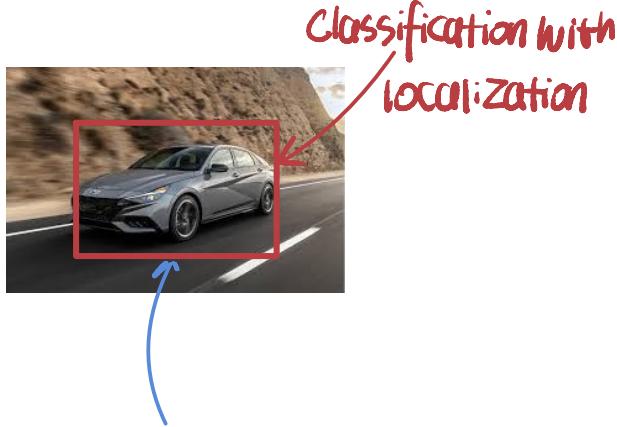
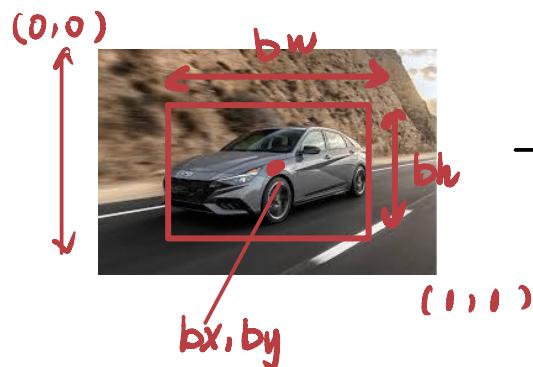
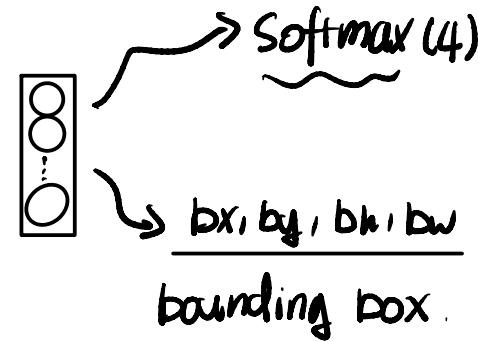


Image classification:  
"Car"

classification with localization



→ ConvNet →



Defining the target label  $y$

"Don't care"

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix} \leftarrow \text{is there any object?}$$



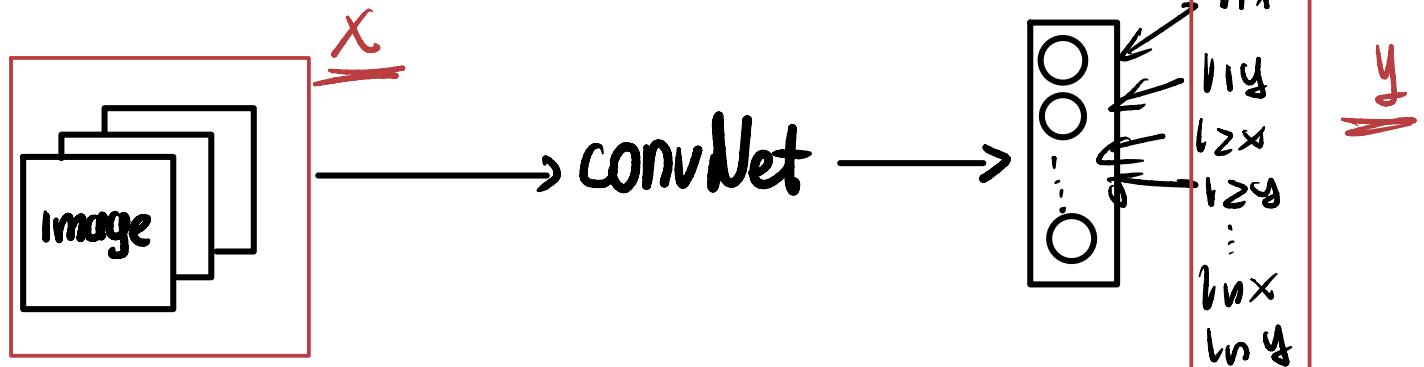
$$\begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ - \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

No object

$$\begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix}$$

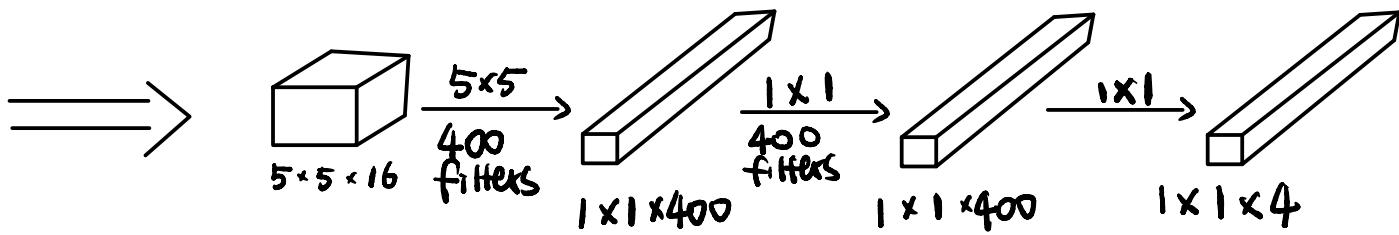
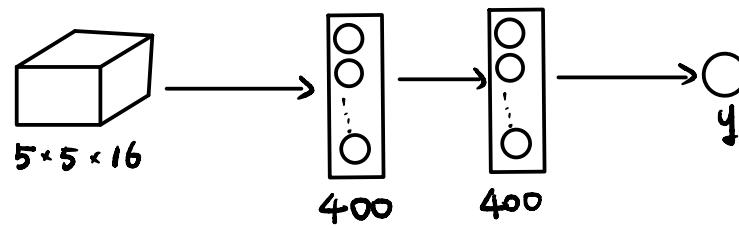
$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_n - y_n)^2 & \text{if } y_i = 1 \\ (\hat{y}_1 - y_1)^2 & \text{if } y_1 = 0 \end{cases}$$

# Landmark Detection

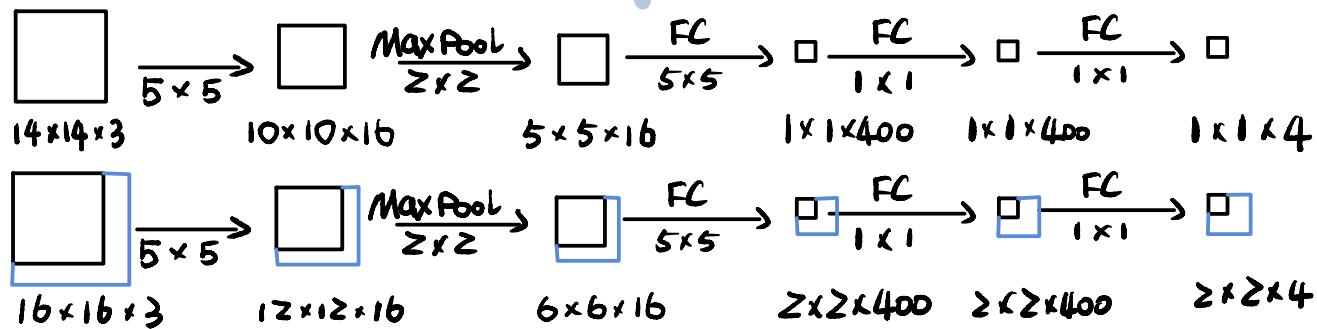


# Object Detection

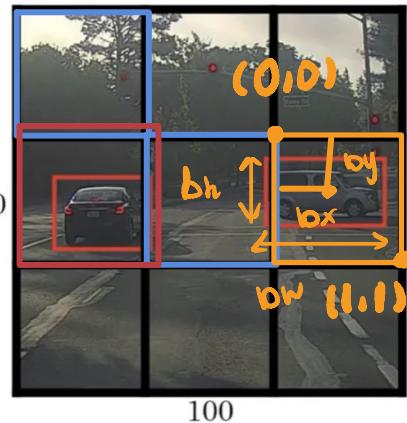
Tuning FC layer into Convolutional layers:



Convolution implementation of Sliding windows



## Bounding Box Predictions



\* Label for training for each grid cell:

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_w \\ b_h \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_w \\ b_h \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$$

} between 0 and 1  
} Could be > 1

can be finer grids.

$\rightarrow x$  (e.g.  $100 \times 100 \times 3$ )  $\rightarrow$  CNN  $\rightarrow$   $3 \times 3 \times 8$   $y$   
input image Target output

## Intersection Over Union:

Intersection over Union (IoU) is used to evaluate object localization

Prediction ( $\hat{y}$ )

IoU =  $\frac{\text{Size of Intersection}}{\text{Size of Union}}$  =  $\frac{\text{Size of } \cap}{\text{Size of } \cup}$

label ( $y$ )

## Non-max Suppression algorithm

- 1) Discard all boxes with  $p_c \leq 0.6$
- 2) While there are any remaining boxes:
  - Pick the box with the largest  $p_c$  output that as a prediction
  - Discard any remaining box with  $\text{IoU} > 0.5$  with the box output in the previous step.

Conduct Non-max Suppression independently for each object categories if the model detects multiple categories.

## Anchor Boxes:

Allowing one grid to contain multiple objects.

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint

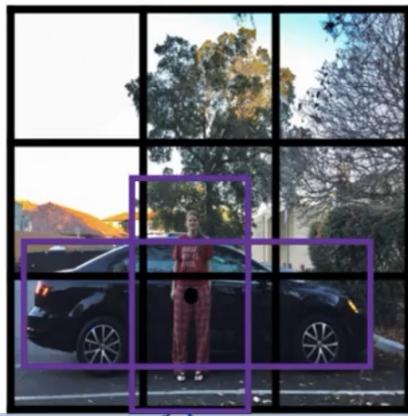
With  $\geq$  anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU

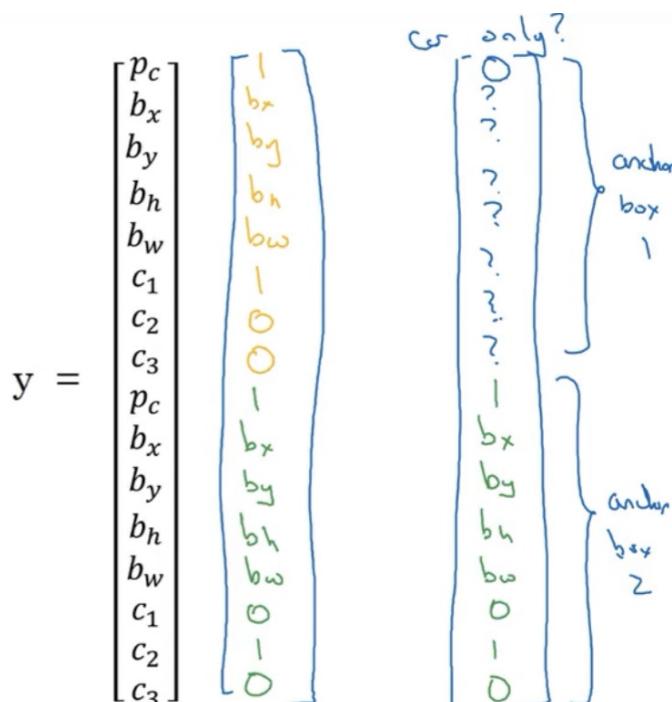
In this case: output  $y = 3 \times 3 \times 16$

$$\sum_{\text{grid cells}} \times 8 \quad \# \text{ anchor boxes}$$

Anchor box example



Anchor box 1: Anchor box 2:



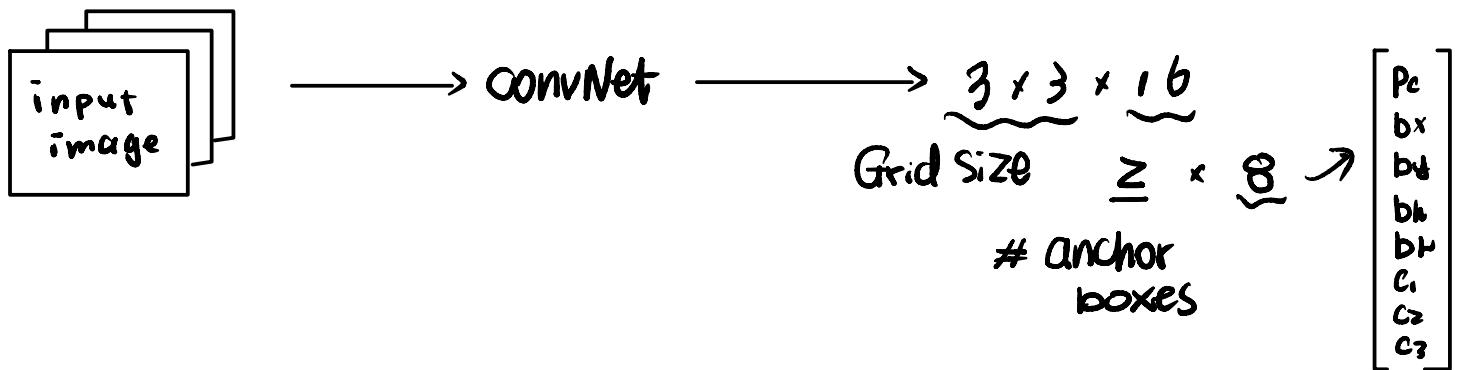
Andrew Ng

The number and shapes of anchor boxes are defined before training, manually or algorithmically

# YOLO Algorithm

Training :

Specify : Object categories (e.g. 1. pedestrian, 2. car, 3. motorcycle).  
anchor boxes



Making Predictions :

image  $\longrightarrow$  trained ConvNet  $\longrightarrow \frac{3 \times 3 \times 2 \times 8}{\text{Predictions}}$

\* Outputting the non-max suppressed outputs :

- For each grid cell, get  $\geq$  predicted bounding boxes
- Get rid of low probability predictions
- For each class, use non-max suppression to generate final predictions.

Region Proposal = idea of semantic segmentation + CNN on segmentations to detect objects and boxes

R-CNN : Propose Regions. Classify proposed regions one at a time  
output label + bounding box

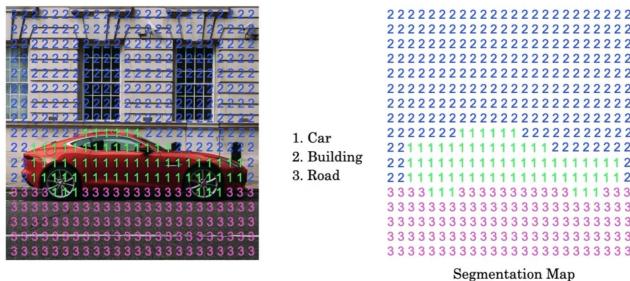
Fast R-CNN : Propose Regions, Use convolutional implementation of sliding window to classify all proposed regions.

Faster R-CNN : Use convolution network to propose regions

# Semantic Segmentation with U-Net

## Semantic Segmentation :

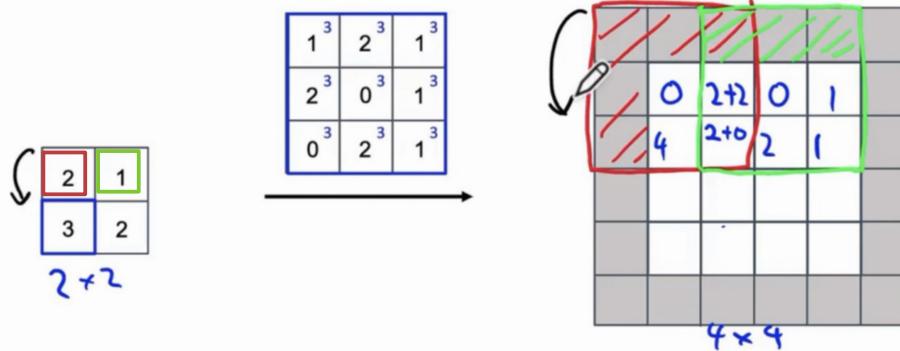
## Per-pixel class labels



Andrew Ng

Challenge: Normal CNN would gradually decrease  $w, h$  and increase  $n_c$ . In order to output per-pixel class labels, we need a type of transformations to bring up  $h, w$  sizes back to the original dimensions.

## → Transpose Convolution



filter  $f \times f = 3 \times 3$  padding  $p = 1$  stride  $s = 2$

Andrew Ng

# U-Net

